

## Lecture 2: 8 January, 2026

Madhavan Mukund

<https://www.cmi.ac.in/~madhavan>

Data Mining and Machine Learning  
January–April 2026

# Market-Basket Analysis

- People who buy  $X$  also tend to buy  $Y$
- Rearrange products on display based on customer patterns

DATA MINING  
VS  
MACHINE LEARNING

# Market-Basket Analysis

- People who buy  $X$  also tend to buy  $Y$
- Rearrange products on display based on customer patterns
  - The diapers and beer legend
  - The true story, <http://www.dssresources.com/newsletters/66.php>

# Market-Basket Analysis

- People who buy  $X$  also tend to buy  $Y$
- Rearrange products on display based on customer patterns
  - The diapers and beer legend
  - The true story, <http://www.dssresources.com/newsletters/66.php>
- Applies in more abstract settings
  - Items are concepts, basket is a set of concepts in which a student does badly
    - Students with difficulties in concept  $A$  also tend to misunderstand concept  $B$
  - Items are words, transactions are documents

# Formal setting

- Set of **items**  $I = \{i_1, i_2, \dots, i_N\}$
- A **transaction** is a set  $t \subseteq I$  of items
- Set of transactions  $T = \{t_1, t_2, \dots, t_M\}$

# Formal setting

- Set of **items**  $I = \{i_1, i_2, \dots, i_N\}$
- A **transaction** is a set  $t \subseteq I$  of items
- Set of transactions  $T = \{t_1, t_2, \dots, t_M\}$
- Identify **association rules**  $X \rightarrow Y$ 
  - $X, Y \subseteq I, X \cap Y = \emptyset$
  - If  $X \subseteq t_j$  then it is likely that  $Y \subseteq t_j$

$X, Y$  are sets

$X \rightarrow \cancel{X \cup Y}$

Actual  
transaction  
has  $X \cup Y$

# Formal setting

- Set of **items**  $I = \{i_1, i_2, \dots, i_N\}$
- A **transaction** is a set  $t \subseteq I$  of items
- Set of transactions  $T = \{t_1, t_2, \dots, t_M\}$
- Identify **association rules**  $X \rightarrow Y$ 
  - $X, Y \subseteq I, X \cap Y = \emptyset$
  - If  $X \subseteq t_j$  then it is likely that  $Y \subseteq t_j$
- Two thresholds
  - How frequently does  $X \subseteq t_j$  imply  $Y \subseteq t_j$ ? **Likelihood**
  - How significant is this pattern overall?

# Setting thresholds

- For  $Z \subseteq I$ ,  $Z.\text{count} = |\{t_j \mid Z \subseteq t_j\}|$

Number of transactions containing  $Z$



# Setting thresholds

- For  $Z \subseteq I$ ,  $Z.\text{count} = |\{t_j \mid Z \subseteq t_j\}|$
  - How frequently does  $X \subseteq t_j$  imply  $Y \subseteq t_j$ ?
    - Fix a confidence level  $\chi$
    - Want  $\frac{(X \cup Y).\text{count}}{X.\text{count}} \geq \chi$
- ≤ 1

# Setting thresholds

- For  $Z \subseteq I$ ,  $Z.\text{count} = |\{t_j \mid Z \subseteq t_j\}|$
- How frequently does  $X \subseteq t_j$  imply  $Y \subseteq t_j$ ?
  - Fix a **confidence level**  $\chi$
  - Want  $\frac{(X \cup Y).\text{count}}{X.\text{count}} \geq \chi$
- How significant is this pattern overall?
  - Fix a **support level**  $\sigma$
  - Want  $\frac{(X \cup Y).\text{count}}{M} \geq \sigma$

total no of transactions

# Setting thresholds

- For  $Z \subseteq I$ ,  $Z.\text{count} = |\{t_j \mid Z \subseteq t_j\}|$
- How frequently does  $X \subseteq t_j$  imply  $Y \subseteq t_j$ ?
  - Fix a **confidence level**  $\chi$
  - Want  $\frac{(X \cup Y).\text{count}}{X.\text{count}} \geq \chi$
- How significant is this pattern overall?
  - Fix a **support level**  $\sigma$
  - Want  $\frac{(X \cup Y).\text{count}}{M} \geq \sigma$
- Given sets of items  $I$  and transactions  $T$ , with confidence  $\chi$  and support  $\sigma$ , find all valid association rules  $X \rightarrow Y$

$$I = \{i_1, \dots, i_n\}$$

$$T = \{t_1, \dots, t_m\}$$

$$X.\text{count} \leq \frac{(X \cup Y).\text{count}}{\chi}$$

# Frequent itemsets

- $X \rightarrow Y$  is interesting only if  $(X \cup Y).count \geq \sigma \cdot M$

- First identify all frequent itemsets

- $Z \subseteq I$  such that  $Z.count \geq \sigma \cdot M$

set of items

$$\frac{Z.count}{M} \geq \sigma$$

# Frequent itemsets

- $X \rightarrow Y$  is interesting only if  $(X \cup Y).count \geq \sigma \cdot M$
- First identify all frequent itemsets
  - $Z \subseteq I$  such that  $Z.count \geq \sigma \cdot M$
- Naïve strategy: maintain a counter for each  $Z$ 
  - For each  $t_j \in T$ 
    - For each  $Z \subseteq t_j$
    - Increment the counter for  $Z$
  - After scanning all transactions, keep  $Z$  with  $Z.count \geq \sigma \cdot M$

Dictionary

$$Z \subseteq I$$

$10^7$  items

$2^{10^7}$  subsets

# Frequent itemsets

- $X \rightarrow Y$  is interesting only if  $(X \cup Y).count \geq \sigma \cdot M$
- First identify all **frequent itemsets**
  - $Z \subseteq I$  such that  $Z.count \geq \sigma \cdot M$
- Naïve strategy: maintain a counter for each  $Z$ 
  - For each  $t_j \in T$   
For each  $Z \subseteq t_j$   
Increment the counter for  $Z$
  - After scanning all transactions, keep  $Z$  with  $Z.count \geq \sigma \cdot M$
- Need to maintain  $2^{|I|}$  counters
  - Infeasible amount of memory
  - Can we do better?

# Sample calculation

- Let's assume a bound on each  $t_i \in \mathcal{T}$ 
  - No transaction has more than 10 items
- Say  $N = |I| = 10^6$ ,  $M = |\mathcal{T}| = 10^9$ ,  $\sigma = 0.01$ 
  - Number of possible subsets to count is  $\sum_{i=1}^{10} \binom{10^6}{i}$

# Sample calculation

- Let's assume a bound on each  $t_i \in T$ 
  - No transaction has more than 10 items
- Say  $N = |I| = 10^6$ ,  $M = |T| = 10^9$ ,  $\sigma = 0.01$ 
  - Number of possible subsets to count is  $\sum_{i=1}^{10} \binom{10^6}{i}$
- A singleton subset that is frequent is an item that appears in at least  $10^7$  transactions

$\{i_{72}\}$   $[i_{329}]$

$$0.01 \times 10^9 = 10^7$$

	1	2	...	10
$t_1$	$i_{72}$			
				$i_{72}$
			$i_{72}$	
$\vdots$				
$t_{10^9}$				

$10^9 \times 10 = 10^{10}$

$$\frac{10^{10}}{10^7} = 1000$$



# Sample calculation

- Let's assume a bound on each  $t_i \in \mathcal{T}$ 
  - No transaction has more than 10 items
- Say  $N = |I| = 10^6$ ,  $M = |\mathcal{T}| = 10^9$ ,  $\sigma = 0.01$ 
  - Number of possible subsets to count is  $\sum_{i=1}^{10} \binom{10^6}{i}$
- A singleton subset that is frequent is an item that appears in at least  $10^7$  transactions
- Totally,  $\mathcal{T}$  contains at most  $10^{10}$  items
- At most  $10^{10}/10^7 = 1000$  items are frequent!
- How can we exploit this?

- Clearly, if  $Z$  is frequent, so is every subset  $Y \subseteq Z$

# Apriori

- Clearly, if  $Z$  is frequent, so is every subset  $Y \subseteq Z$
- We exploit the contrapositive

## Apriori observation

If  $Z$  is not a frequent itemset, no superset  $Y \supseteq Z$  can be frequent

- Clearly, if  $Z$  is frequent, so is every subset  $Y \subseteq Z$
- We exploit the contrapositive

## Apriori observation

If  $Z$  is not a frequent itemset, no superset  $Y \supseteq Z$  can be frequent

- For instance, in our earlier example, every frequent itemset must be built from the 1000 frequent items
- In particular, for any frequent pair  $\{x, y\}$ , both  $\{x\}$  and  $\{y\}$  must be frequent
- Build frequent itemsets bottom up, size 1, 2, ...

# Apriori algorithm

- $F_i$  : frequent itemsets of size  $i$  — Level  $i$

# Apriori algorithm

- $F_i$  : frequent itemsets of size  $i$  — Level  $i$
- $F_1$ : Scan  $T$ , maintain a counter for each  $x \in I$

# Apriori algorithm

- $F_i$  : frequent itemsets of size  $i$  — Level  $i$
- $F_1$ : Scan  $T$ , maintain a counter for each  $x \in I$
- $C_2 = \{\{x, y\} \mid x, y \in F_1\}$ : Candidates in level 2

$x \neq y$

$10^3 \times 10^3$

# Apriori algorithm

- $F_i$  : frequent itemsets of size  $i$  — Level  $i$
- $F_1$ : Scan  $T$ , maintain a counter for each  $x \in I$
- $C_2 = \{\{x, y\} \mid x, y \in F_1\}$ : Candidates in level 2
- $F_2$ : Scan  $T$ , maintain a counter for each  $X \in C_2$



# Apriori algorithm

- $F_i$  : frequent itemsets of size  $i$  — Level  $i$
- $F_1$ : Scan  $T$ , maintain a counter for each  $x \in I$
- $C_2 = \{\{x, y\} \mid x, y \in F_1\}$ : Candidates in level 2
- $F_2$ : Scan  $T$ , maintain a counter for each  $X \in C_2$
- $C_3 = \{\{x, y, z\} \mid \{x, y\}, \{x, z\}, \{y, z\} \in F_2\} \Rightarrow$
- $F_3$ : Scan  $T$ , maintain a counter for each  $X \in C_3$

Implies  
↓  
 $x \in F_1, y \in F_1, z \in F_1$

# Apriori algorithm

- $F_i$  : frequent itemsets of size  $i$  — Level  $i$
- $F_1$ : Scan  $T$ , maintain a counter for each  $x \in I$
- $C_2 = \{\{x, y\} \mid x, y \in F_1\}$ : Candidates in level 2
- $F_2$ : Scan  $T$ , maintain a counter for each  $X \in C_2$
- $C_3 = \{\{\underline{x}, y, z\} \mid \{x, y\}, \{x, z\}, \{y, z\} \in F_2\}$  — Naively check  $(10^7)^3$  triples
- $F_3$ : Scan  $T$ , maintain a counter for each  $X \in C_3$
- ...
- $C_k$  = subsets of size  $k$ , every  $(k-1)$ -subset is in  $F_{k-1}$
- $F_k$ : Scan  $T$ , maintain a counter for each  $X \in C_k$
- ...

Bottleneck

$$C_2 \approx F_1 \times F_1$$

$$F_2 \times F_2 \times F_2$$

# Apriori algorithm

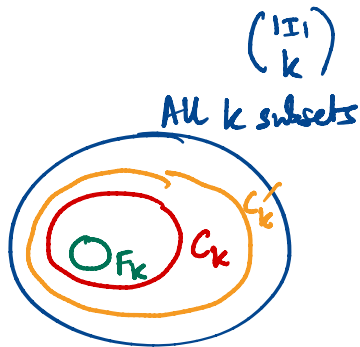
- $C_k$  = subsets of size  $k$ , every  $(k-1)$ -subset is in  $F_{k-1}$
- How do we generate  $C_k$ ?

# Apriori algorithm

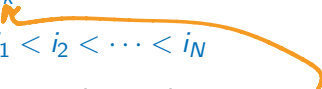
- $C_k$  = subsets of size  $k$ , every  $(k-1)$ -subset is in  $F_{k-1}$
- How do we generate  $C_k$ ?
- Naïve: enumerate subsets of size  $k$  and check each one
  - Expensive!

# Apriori algorithm

- $C_k$  = subsets of size  $k$ , every  $(k-1)$ -subset is in  $F_{k-1}$
- How do we generate  $C_k$ ?
- Naïve: enumerate subsets of size  $k$  and check each one
  - Expensive!
- **Observation:** Any  $C'_k \supseteq C_k$  will do as a candidate set



# Apriori algorithm

- $C_k$  = subsets of size  $k$ , every  $(k-1)$ -subset is in  $F_{k-1}$
  - How do we generate  $C_k$ ?
  - Naïve: enumerate subsets of size  $k$  and check each one
    - Expensive!
  - **Observation:** Any  $C'_k \supseteq C_k$  will do as a candidate set
  - Items are ordered:  $i_1 < i_2 < \dots < i_N$
  - List each itemset in ascending order — canonical representation
- 

# Apriori algorithm

- $C_k$  = subsets of size  $k$ , every  $(k-1)$ -subset is in  $F_{k-1}$
- How do we generate  $C_k$ ?
- Naïve: enumerate subsets of size  $k$  and check each one
  - Expensive!
- **Observation:** Any  $C'_k \supseteq C_k$  will do as a candidate set
- Items are ordered:  $i_1 < i_2 < \dots < i_N$
- List each itemset in ascending order — canonical representation
- Merge two  $(k-1)$ -subsets if they differ in last element

- $X = \{i_1, i_2, \dots, i_{k-2}, i_{k-1}\}$

- $X' = \{i_1, i_2, \dots, i_{k-2}, i'_{k-1}\}$

- $\text{Merge}(X, X') = \{\underline{i_1, i_2, \dots, i_{k-2}}, \underline{i_{k-1}, i'_{k-1}}\}$



Only 2  $k-1$  subsets are guaranteed to be in  $F_{k-1}$

# Apriori algorithm

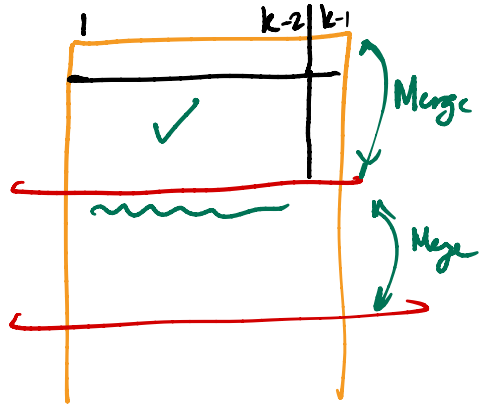
- $\text{Merge}(X, X') = \{i_1, i_2, \dots, i_{k-2}, i_{k-1}, i'_{k-1}\}$ 
  - $X = \{i_1, i_2, \dots, i_{k-2}, i_{k-1}\}$
  - $X' = \{i_1, i_2, \dots, i_{k-2}, i'_{k-1}\}$



# Apriori algorithm

- $\text{Merge}(X, X') = \{i_1, i_2, \dots, i_{k-2}, i_{k-1}, i'_{k-1}\}$ 
  - $X = \{i_1, i_2, \dots, i_{k-2}, i_{k-1}\}$
  - $X' = \{i_1, i_2, \dots, i_{k-2}, i'_{k-1}\}$
- $C'_k = \{\text{Merge}(X, X') \mid X, X' \in F_{k-1}\}$

Write  $F_{k-1}$  in dictionary  
order



# Apriori algorithm

- $\text{Merge}(X, X') = \{i_1, i_2, \dots, i_{k-2}, i_{k-1}, i'_{k-1}\}$ 
  - $X = \{i_1, i_2, \dots, i_{k-2}, i_{k-1}\}$
  - $X' = \{i_1, i_2, \dots, i_{k-2}, i'_{k-1}\}$
- $C'_k = \{\text{Merge}(X, X') \mid X, X' \in F_{k-1}\}$
- **Claim**  $C_k \subseteq C'_k$ 
  - Suppose  $Y = \{i_1, i_2, \dots, i_{k-1}, i_k\} \in C_k$
  - $X = \{i_1, i_2, \dots, i_{k-2}, i_{k-1}\} \in F_{k-1}$  and  $X' = \{i_1, i_2, \dots, i_{k-2}, i_k\} \in F_{k-1}$
  - $Y = \text{Merge}(X, X') \in C'_k$



# Apriori algorithm

- $\text{Merge}(X, X') = \{i_1, i_2, \dots, i_{k-2}, i_{k-1}, i'_{k-1}\}$ 
  - $X = \{i_1, i_2, \dots, i_{k-2}, i_{k-1}\}$
  - $X' = \{i_1, i_2, \dots, i_{k-2}, i'_{k-1}\}$
- $C'_k = \{\text{Merge}(X, X') \mid X, X' \in F_{k-1}\}$
- **Claim**  $C_k \subseteq C'_k$ 
  - Suppose  $Y = \{i_1, i_2, \dots, i_{k-1}, i_k\} \in C_k$
  - $X = \{i_1, i_2, \dots, i_{k-2}, i_{k-1}\} \in F_{k-1}$  and  $X' = \{i_1, i_2, \dots, i_{k-2}, i_k\} \in F_{k-1}$
  - $Y = \text{Merge}(X, X') \in C'_k$
- Can generate  $C'_k$  efficiently
  - Arrange  $F_{k-1}$  in dictionary order
  - Split into blocks that differ on last element
  - Merge all pairs within each block

# Apriori algorithm

- $C_1 = \{\{x\} \mid x \in I\}$
- $F_1 = \{Z \mid Z \in C_1, Z.\text{count} \geq \sigma \cdot M\}$
- For  $k \in \{2, 3, \dots\}$ 
  - $C'_k = \{\text{Merge}(X, X') \mid X, X' \in F_{k-1}\}$
  - $F_k = \{Z \mid Z \in C'_k, Z.\text{count} \geq \sigma \cdot M\}$

$F_1$      ~~$\{x\}$~~   
           ~~$\{y\}$~~   
           $\{x, y\}$

# Apriori algorithm

- $C_1 = \{\{x\} \mid x \in I\}$
- $F_1 = \{Z \mid Z \in C_1, Z.\text{count} \geq \sigma \cdot M\}$
- For  $k \in \{2, 3, \dots\}$ 
  - $C'_k = \{\text{Merge}(X, X') \mid X, X' \in F_{k-1}\}$
  - $F_k = \{Z \mid Z \in C'_k, Z.\text{count} \geq \sigma \cdot M\}$
- When do we stop?

# Apriori algorithm

- $C_1 = \{\{x\} \mid x \in I\}$
- $F_1 = \{Z \mid Z \in C_1, Z.\text{count} \geq \sigma \cdot M\}$
- For  $k \in \{2, 3, \dots\}$ 
  - $C'_k = \{\text{Merge}(X, X') \mid X, X' \in F_{k-1}\}$
  - $F_k = \{Z \mid Z \in C'_k, Z.\text{count} \geq \sigma \cdot M\}$
- When do we stop?
- $k$  exceeds the size of the largest transaction
- $F_k$  is empty

# Apriori algorithm

- $C_1 = \{\{x\} \mid x \in I\}$
- $F_1 = \{Z \mid Z \in C_1, Z.\text{count} \geq \sigma \cdot M\}$
- For  $k \in \{2, 3, \dots\}$ 
  - $C'_k = \{\text{Merge}(X, X') \mid X, X' \in F_{k-1}\}$
  - $F_k = \{Z \mid Z \in C'_k, Z.\text{count} \geq \sigma \cdot M\}$
- When do we stop?
- $k$  exceeds the size of the largest transaction
- $F_k$  is empty

"Exact" computation

Next step: From frequent itemsets to association rules